

# Mapping design values over Canada

Chao Li et al.

March 5, 2018

## Data formatting

We represent the inferred design values (e.g., 100-year extreme rainfall) in a region from the large ensemble CanRCM4 simulations by an  $n \times p$  data matrix  $\mathbf{X}$ , with  $n$  being the ensemble size of the simulations and  $p$  the number of grid cells in the region. The region can be the entire land of Canada or a subregion. We center the data matrix by subtracting the ensemble mean of the design values at each grid cell, that is,  $\mathbf{X}' = (\mathbf{I}_n - \frac{1}{n}\mathbf{1}_n)\mathbf{X}$ , where  $\mathbf{I}_n$  is the identity matrix of size  $n$  and  $\mathbf{1}_n$  is the  $n \times n$  all-ones matrix. To minimize the influence of the varying grid cell areas on the structure of the computed EOFs, we weight each value in  $\mathbf{X}'$  by the grid cell area, that is,  $\mathbf{X}'_w = \mathbf{X}'\mathbf{W}$ , with  $\mathbf{W} = \text{diag}(f_1, f_2, \dots, f_p)$  being the weighting matrix and  $f_i$  the fractional area of the grid cell  $i$ . In the following and unless otherwise stated, we drop the super- and subscripts and denote the centered and weighted design value matrix by  $\mathbf{X}$ . The spatial covariance of  $\mathbf{X}$  is estimated by  $\mathbf{C} = \mathbf{X}^T \mathbf{X}$ .

## EOF with data matrix

EOF analysis can be computed by singular value decomposition (SVD) of the data matrix  $\mathbf{X}$ , or by eigendecomposition of the covariance matrix  $\mathbf{C}$ . According

to SVD, the  $n \times p$  matrix  $\mathbf{X}$  can be factored as

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \quad (1)$$

where  $\mathbf{U}$  is an  $n \times n$  orthogonal matrix whose columns  $\mathbf{u}_i$  ( $i = 1, 2, \dots, n$ ) are the EOF coefficients which are often as well referred to as principal components (PCs);  $\mathbf{V}$  is a  $p \times p$  orthogonal matrix whose columns  $\mathbf{v}_i$  ( $i = 1, 2, \dots, p$ ) provides the EOFs; and  $\mathbf{\Sigma}$  is an  $n \times p$  diagonal matrix of the form  $\text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r, 0, \dots, 0)$  with  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  and  $r = \text{rank}(\mathbf{X}) \leq \min(n, p)$ . In the above,  $\sigma_1^2, \dots, \sigma_r^2$  are the eigenvalues of  $\mathbf{X}^T \mathbf{X}$  (also of the covariance matrix  $\mathbf{C}$ ).

## EOF with covariance matrix

The covariance matrix  $\mathbf{C}$  (of size  $p \times p$ ) has an eigendecomposition in the form of

$$\mathbf{C} = \mathbf{V}\mathbf{A}\mathbf{V}^T \quad (2)$$

where  $\mathbf{V}$  is a  $p \times p$  orthogonal matrix containing the EOFs  $\mathbf{v}_i$  ( $i = 1, 2, \dots, p$ ) in its columns;  $\mathbf{A}$  is a  $p \times p$  diagonal matrix containing the eigenvalues of  $\mathbf{C}$  down the diagonal, that is,  $\mathbf{A} = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2, 0, \dots, 0)$  with  $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_r^2 > 0$  and  $r \leq \min(n, p)$ . Projecting  $\mathbf{X}$  onto the EOFs by  $\mathbf{A} = \mathbf{X}\mathbf{V}$  obtains, in the columns of  $\mathbf{A}$ , the EOF coefficients. Unlike those computed via SVD, here the EOF coefficients  $\alpha_i$  ( $i = 1, 2, \dots, p$ ) are not necessarily unitary. With the EOFs and their coefficients, the data matrix  $\mathbf{X}$  can be factored as

$$\mathbf{X} = \mathbf{A}\mathbf{V}^{-1} = \mathbf{A}\mathbf{V}^T \quad (3)$$

## EOF reconstruction: the basis of the mapping approach

The EOFs describe the spatial dependence structure of the design value field, while the EOF coefficients reflect the variation of the EOF modes due to internal variability. The eigenvalues give a measure of the importance of each EOF mode. Only the first  $r$  EOF modes with nonzero eigenvalues carry information. Neglecting the redundant ones, the data matrix  $\mathbf{X}$  can be rewritten in vector form from (1) as follows

$$\mathbf{X} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \quad (4)$$

Likewise, from (3), we have

$$\mathbf{X} = \sum_{i=1}^r \alpha_i \mathbf{v}_i^T \quad (5)$$

The reconstructions (4-5) state that the design value field can be fully represented by a set of basis functions of space (i.e., EOFs) and a set of amplitude functions of internal variability (i.e., EOF coefficients). This provides the fundamental of our approach to design value mapping.

In fact, it is a few leading EOFs that are capturing large-scale and low-frequency features of the field, with the remaining mostly representing random noise. The reconstruction is therefore typically implemented by truncating the above sums at some  $m$  ( $\ll r$ ) such that only the important large-scale features are retained while the random noise is filtered. In practice, the truncation level  $m$  is determined by prescribing a certain amount of the explained variance (e.g., 95%) and choosing the first  $m$  EOFs that explain altogether this amount of variance. The explained variance by the first  $m$  EOFs is computed as a percent by

$$\frac{100 \sum_{i=1}^m \sigma_i^2}{\sum_{i=1}^r \sigma_i^2} \%$$

## The design value mapping method

Given the relatively sparse distribution of meteorological observing stations in Canada plus the fact that many of these stations cover too short periods to allow reliable estimates of design values, traditional interpolation methods such as spline interpolation and Kriging become impractical. We thus base the mapping method on EOF reconstruction. To be specific, the method obtains leading EOF modes of design values from a 50-member ensemble of CanRCM4 simulations, constrains the amplitudes of these modes (i.e., EOF coefficients) by the available observed design values via a least squares regression, and then combines them together to reconstruct design values at unobserved locations. The underlying assumption is that CanRCM4 simulations can reasonably resolve the overall spatial structure of observed design value field, despite the possibility that the magnitudes of design values might differ between models and observations.

The mapping method needs a dataset of observed design values gridded consistently with CanRCM4 simulations. For that, we grid the design values estimated from observations at stations with sufficient data by averaging available values within each CanRCM4 grid cell ( $\sim 50 \text{ km} \times 50 \text{ km}$ ). Grid cells with no observations are flagged as missing. The gridded observations are then organized into a matrix  $\mathbf{y}_g$  (which is actually a  $1 \times p$  vector) and weighted by grid cell area, both following the way as for CanRCM4 simulations. We use the subscript “ $g$ ” to make it explicit that the matrix contain gaps (i.e., missing values). For operational convenience, we treat  $\mathbf{y}_g$  as a *column* vector in the following.

With  $\mathbf{X}$  and  $\mathbf{y}_g$ , the mapping method is implemented as follows:

- 1) Calculate EOFs  $\mathbf{V}$  and eigenvalues, by SVD of  $\mathbf{X}$  or eigendecomposition of  $\mathbf{X}^T \mathbf{X}$ .

2) For a given number  $k$  of EOFs, find the observation-constrained EOF coefficients  $\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)^T$  via a least squares regression,  $\mathbf{y}_{g,\perp} = \mathbf{V}_{1:k,\perp}\beta + \epsilon$ , where subscript  $1:k$  means that the first  $k$  columns of  $\mathbf{V}$  (which correspond to the first  $k$  EOFs) are involved in the regression; and the subscript  $\perp$  indicates collapsed versions of the corresponding matrices with rows that do not have observations removed. The solution  $\hat{\beta}$  is found by minimizing the objective function  $\psi_k = \|\mathbf{y}_{g,\perp} - \mathbf{V}_{1:k,\perp}\beta\|^2$ , giving  $\hat{\beta} = (\mathbf{V}_{1:k,\perp}^T \mathbf{V}_{1:k,\perp})^{-1} \mathbf{V}_{1:k,\perp}^T \mathbf{y}_{g,\perp}$ .

3) Repeat 2) by increasing  $k$  to  $k+1$  until the change of the objective function is  $< 10^{-3}$  or the explained variance by the involved EOFs is  $\geq 95\%$ . If neither of the criteria is fulfilled but the number of involved EOFs becomes greater than 10% of the number of the grid cells having observations, the method is marked as not converged, otherwise denote the convergence step as  $l$  and proceed to 4).

4) Reconstruct the full design value field and remove the effect of areal weighting from the reconstructions by conducting  $\mathbf{y} = \mathbf{V}_{1:l}\hat{\beta}\mathbf{W}^{-1}$ . In order to respect grid cells with observations, relevant reconstructed values are substituted by the corresponding observations. The reconstructions are then used to plot a map for design values with full coverage over Canada.